

# Lecture 5: Model-Free Control

David Silver

All previous lectures lead to this lecture.

You drop your robot/agent into an unknown environment, we know nothing about how the environment works, how can we max the rewards? The core techniques we used in this lecture is built on lec4.

In future lectures, we will talk about how scale up.

# Outline

- 1 Introduction
- 2 On-Policy Monte-Carlo Control
- 3 On-Policy Temporal-Difference Learning
- 4 Off-Policy Learning
- 5 Summary

# Model-Free Reinforcement Learning

lec 3:

Planning (prediction, control) by DP.  
Solve a known MDP.

Lec 4 ■ Last lecture:

Drop your agent in an unknown MDP with a given policy, how to evaluate this policy, how much rewards we can get if following the behaviors of this policy.

- **Model-free prediction**
- **Estimate** the value function of an **unknown** MDP

Lec 5 ■ This lecture:

- **Model-free control**
- **Optimise** the value function of an **unknown** MDP

Find  $v_*$ ,  $q_*$

We use same tools, we iterate them and find the best possible behaviors.

# Uses of Model-Free Control

Why interesting? Why useful?

Some **example problems** that can be modelled as MDPs

- Elevator
- Parallel Parking
- Ship Steering
- Bioreactor
- Helicopter
- Aeroplane Logistics
- Robocup Soccer
- Quake
- Portfolio management
- Protein Folding
- Robot walking
- Game of Go

**For most of these problems, either:**

These problems are unknown to use. We don't know the environment so have to use model-free MDPs.

- **MDP model is unknown, but experience can be sampled**
- **MDP model is known, but is too big to use, except by samples**

**Model-free control** can solve these problems

# On and Off-Policy Learning

## ■ On-policy learning

Follow the behaviors we learn from this job.

You get a policy, you follow that policy. While following it, you learn about that policy.

- “Learn on the job”

- Learn about policy  $\pi$  from experience sampled from  $\pi$

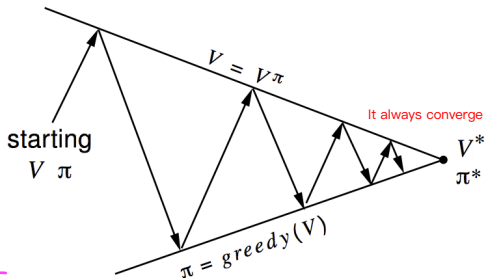
## ■ Off-policy learning

- “Look over someone’s shoulder” Follow someone else’s behaviors.

- Learn about policy  $\pi$  from experience sampled from  $\mu$

The robot/agent can learn not only from itself’s experience but also others’. Other can be other robot/agent or even human demonstrations.

# Generalised Policy Iteration (Refresher)



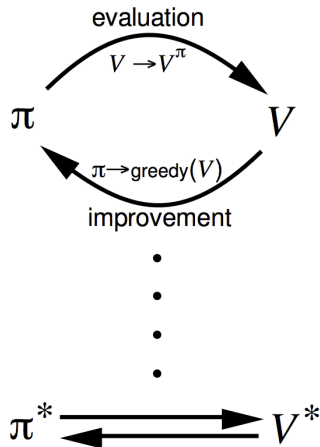
Firstly

**Policy evaluation** Estimate  $v_\pi$

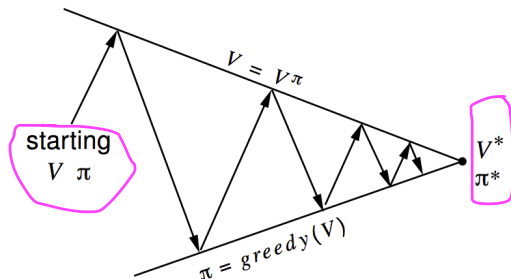
e.g. Iterative policy evaluation

**Policy improvement** Generate  $\pi' \geq \pi$

e.g. Greedy policy improvement



# Generalised Policy Iteration With Monte-Carlo Evaluation



Policy evaluation Monte-Carlo policy evaluation,  $V = v_\pi$

Take mean value

Policy improvement Greedy policy improvement?

Will this MC + Greedy combination work? It has 2 issues:

1. Evaluation Step: If we use status value function  $V$ , we need look ahead one step to use value func of next state while need to know the action to be taken but this makes it not model-free (because asking for know the action). Solve: Use Action value function instead (see next slide). One more small issue: It will be slow. It needs lots of efforts to do so. This can be improved by TD later.
2. Improvement Step: If you always greedy, you will not explore the whole state space. So there might be some potential you never see.

# Model-Free Policy Iteration Using Action-Value Function

To replace State-Value function  $V$

- Greedy policy improvement over  $V(s)$  **requires** model of MDP

$$\pi'(s) = \operatorname{argmax}_{a \in \mathcal{A}} \mathcal{R}_s^a + \mathcal{P}_{ss'}^a V(s')$$

Which makes it not model  
free anymore

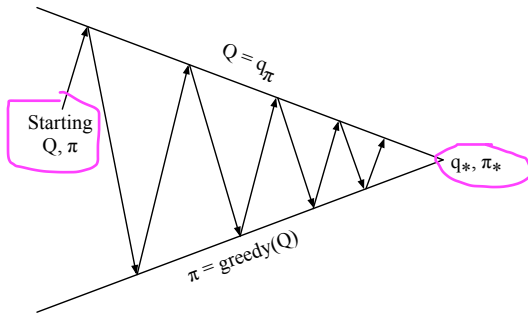
- Greedy policy improvement over  $Q(s, a)$  is **model-free**

$$\pi'(s) = \operatorname{argmax}_{a \in \mathcal{A}} Q(s, a)$$

$Q$  tells us how good to  
take each action  
Then we find the Max. No  
need model here.



# Generalised Policy Iteration with Action-Value Function



Solution to Issue on Policy evaluation step:

Replace  
 $V = v_{\pi}$

Policy evaluation Monte-Carlo policy evaluation,  $Q = q_\pi$

Policy improvement Greedy policy improvement?

How to solve the issue on Improvement step?

# Example of Greedy Action Selection



"Behind one door is tenure - behind the other is flipping burgers at McDonald's."

So if greedy you get stuck with "right door". You never know what is the reward of "2nd/3rd/4th .. time choosing left door". (This is actually the drawbacks of greedy algorithm.)

- There are two doors in front of you.
- Time
- 1 ■ You open the left door and get reward 0  
 $V(\text{left}) = 0$
- 2 ■ You open the right door and get reward +1  
 $V(\text{right}) = +1$  MC: choose right bc Mean=1  
Greedy: choose right bc CurrentReward=1
- 3 ■ You open the right door and get reward +3  
 $V(\text{right}) = +2$  MC: choose right bc Mean=1.5  
Greedy: choose right bc CurrentReward=2
- 4 ■ You open the right door and get reward +2  
 $V(\text{right}) = +2$  MC: choose right bc Mean=5/3  
Greedy: choose right bc CurrentReward=2
- ⋮
- Are you sure you've chosen the best door?

# $\epsilon$ -Greedy Exploration

How to guarantee you visit all states or all actions?

- Simplest idea for ensuring continual exploration But work well and efficiently
- All  $m$  actions are tried with non-zero probability
- With probability  $1 - \epsilon$  choose the greedy action
- With probability  $\epsilon$  choose an action at random

$a = \text{greedy}(Q)$  if  $1 - \epsilon$

$a = \text{random } a$  if  $\epsilon$

$$\pi(a|s) = \begin{cases} \epsilon/m + 1 - \epsilon & \text{if } a^* = \underset{a \in \mathcal{A}}{\operatorname{argmax}} Q(s, a) \text{ greedy} \\ \epsilon/m & \text{otherwise random} \end{cases}$$

Fancy way to show same thing:

$$\pi(a|s) = P[A_t = a | S_t = s]$$

e.g.: if we have 10 action options ( $m=10$ ), named  $a_1, a_2 \dots a_{10}$ , and  $a_{10}$  is best option, and  $\epsilon=30\%$  for random.

Then possibilities for all 10 action options:

$P(a_1)=P(a_2)=P(a_3)=\dots=P(a_9)=0.03$ . They only will be chosen in 30% case.

$P(a_{10})=0.73$ . It will be chosen in 30% case and 70% case.

# $\epsilon$ -Greedy Policy Improvement

## Theorem

For any  $\epsilon$ -greedy policy  $\pi$ , the  $\epsilon$ -greedy policy  $\pi'$  with respect to  $q_\pi$  is an improvement,  $v_{\pi'}(s) \geq v_\pi(s)$

action value func  $q_\pi(s, a) = E_\pi[G_t | S_t = s, A_t = a]$

Take one  
step ahead  
following  
new policy  
 $\pi'$

$$q_\pi(s, \pi'(s)) = \sum_{a \in \mathcal{A}} \pi'(a|s) q_\pi(s, a)$$

$$= \epsilon/m \sum_{a \in \mathcal{A}} q_\pi(s, a) + (1 - \epsilon) \max_{a \in \mathcal{A}} q_\pi(s, a)$$

bc max  $\geq$  mean

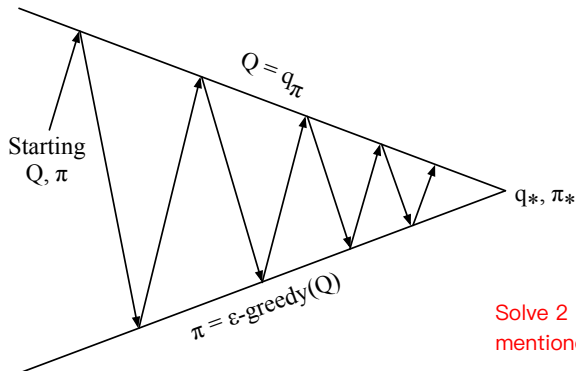
$$\geq \epsilon/m \sum_{a \in \mathcal{A}} q_\pi(s, a) + (1 - \epsilon) \sum_{a \in \mathcal{A}} \frac{\pi(a|s) - \epsilon/m}{1 - \epsilon} q_\pi(s, a)$$

$$= \sum_{a \in \mathcal{A}} \pi(a|s) q_\pi(s, a) = v_\pi(s)$$

Better than taking one step ahead following old policy  $\pi$

Therefore from policy improvement theorem,  $v_{\pi'}(s) \geq v_\pi(s)$

# Monte-Carlo Policy Iteration



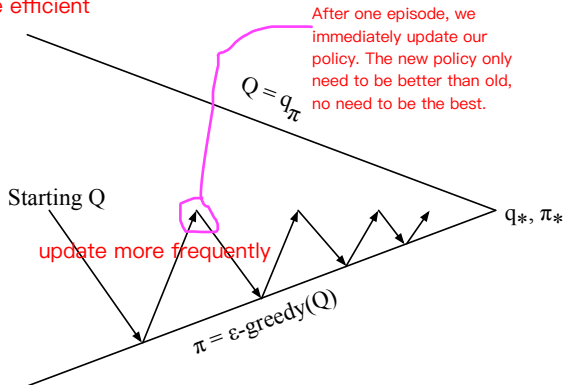
Solve 2 issues we mentioned before.

Policy evaluation Monte-Carlo policy evaluation,  $Q = q_\pi$  Replace  $V=v_\pi$

Policy improvement  $\epsilon\text{-greedy}$  policy improvement Replace greedy.

# Monte-Carlo Control

more efficient



Idea is always act greedy to wrt the most freshest most recent estimated action-value function. After one episode, you can update the value function slightly better, instead of using old estimated action-value function to generate action, use your new updated estimated action-value function to generate behavior.

Every episode:

Policy evaluation Monte-Carlo policy evaluation,  $Q \approx q_\pi$

Policy improvement  $\epsilon$ -greedy policy improvement

How can we guarantee we find the  $\pi_*$ ? We should balance two things. 1 we keep exploring and don't exclude anything which can make it better. 2 asymptotically we get to a policy we're not exploring at all anymore bc the best policy don't include the random behavior.

# GLIE

## Definition

### *Greedy in the Limit with Infinite Exploration (GLIE)*

- All state-action pairs are explored infinitely many times,  
Every action from one state will be tried

$$\lim_{k \rightarrow \infty} N_k(s, a) = \infty$$

- The policy converges on a greedy policy,  
It needs to meet the bellman optimality equation which has a max.

$$\lim_{k \rightarrow \infty} \pi_k(a|s) = \mathbf{1}(a = \operatorname{argmax}_{a' \in \mathcal{A}} Q_k(s, a'))$$

- For example,  $\epsilon$ -greedy is GLIE if  $\epsilon$  reduces to zero at  $\epsilon_k = \frac{1}{k}$

# GLIE Monte-Carlo Control

- Sample  $k$ th episode using  $\pi$ :  $\{S_1, A_1, R_2, \dots, S_T\} \sim \pi$
- For each state  $S_t$  and action  $A_t$  in the episode,

counter:  $N(S_t, A_t) \leftarrow N(S_t, A_t) + 1$

Update  
the mean:  $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{1}{N(S_t, A_t)} (G_t - Q(S_t, A_t))$

- Improve policy based on new action-value function

$$\epsilon \leftarrow 1/k$$

$$\pi \leftarrow \epsilon\text{-greedy}(Q)$$

## Theorem

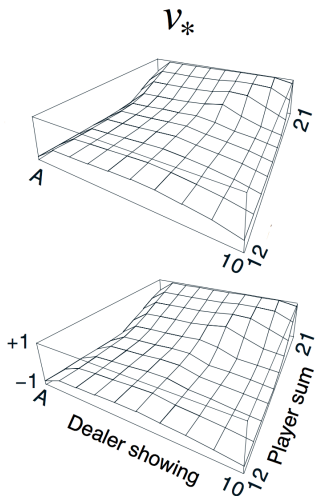
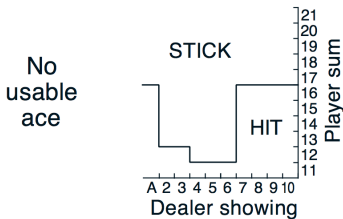
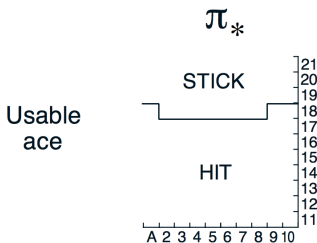
*GLIE Monte-Carlo control converges to the optimal action-value function,  $Q(s, a) \rightarrow q_*(s, a)$*



## Back to the Blackjack Example



# Monte-Carlo Control in Blackjack

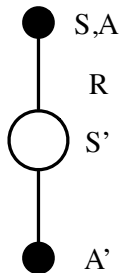


# MC vs. TD Control

- Temporal-difference (TD) learning has several **advantages** over Monte-Carlo (MC)
  - Lower variance
  - Online
  - Incomplete sequences
- Natural idea: use TD instead of MC in our control loop
  - Apply TD to  $Q(S, A)$
  - Use  $\epsilon$ -greedy policy improvement
  - Update every time-step from every episode

# Updating Action-Value Functions with Sarsa

I'm in state  $S$ , I'm considering if I actually take an action, how many rewards I will get, and the value of next action I would take. This is estimated value of that policy and used to update the value of state-action pair I started in.



State-action pair. Start with state  $S$ , choosing an action,

sample from environment to get reward  $R$ ,

end with new state  $S'$ . Then sample our own policy at next state  $S'$ .

TD error

TD target

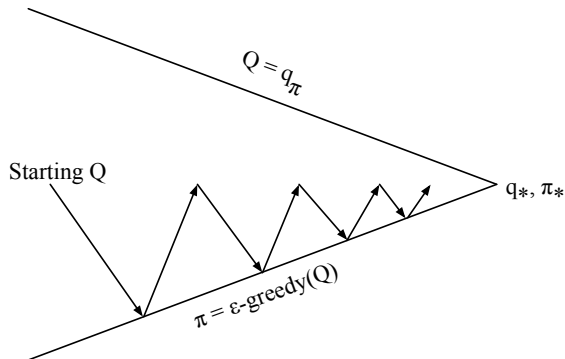
$$Q(S, A) \leftarrow Q(S, A) + \alpha (R + \gamma Q(S', A') - Q(S, A))$$

TD target:

Immediate reward +  
discounted value of next  
state

–  $q$  value where we started

# On-Policy Control With Sarsa



Every time-step:

Policy evaluation Sarsa,  $Q \approx q_\pi$

Policy improvement  $\epsilon$ -greedy policy improvement

# Sarsa Algorithm for On-Policy Control

A lookup table

Initialize  $Q(s, a), \forall s \in \mathcal{S}, a \in \mathcal{A}(s)$ , arbitrarily, and  $Q(\text{terminal-state}, \cdot) = 0$

Repeat (for each episode):

Initialize  $S$

Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)

Repeat (for each step of episode):

Take action  $A$ , observe  $R, S'$  Reward and next state it end up in.

Choose  $A'$  from  $S'$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)

$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma Q(S', A') - Q(S, A)]$

$S \leftarrow S'; A \leftarrow A';$

until  $S$  is terminal

$A'$  is selected  
using current  
policy.

Have a  
new  
policy  
now

# Convergence of Sarsa

## Theorem

*Sarsa converges to the optimal action-value function,* As GLIE-MC

$Q(s, a) \rightarrow q_*(s, a)$ , under the following conditions:

- *GLIE sequence of policies  $\pi_t(a|s)$*
- *Robbins-Monro sequence of step-sizes  $\alpha_t$*

$$\sum_{t=1}^{\infty} \alpha_t = \infty$$

Step size is sufficiently large so you can move your q value as far as you want.

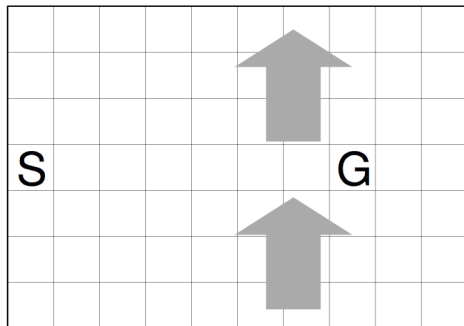
$$\sum_{t=1}^{\infty} \alpha_t^2 < \infty$$

Eventually the changes of your q value becomes smaller and smaller and to 0.

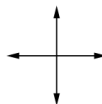
Empirical result: we often don't worry about these 2 so sarsa typically works.

# Windy Gridworld Example

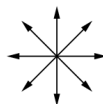
Move from start cell S to goal cell G. Use king's moves. Each move, the wind will move us up by wind strength pieces of cells.



Wind strength: 0 0 0 1 1 1 2 2 1 0



standard  
moves

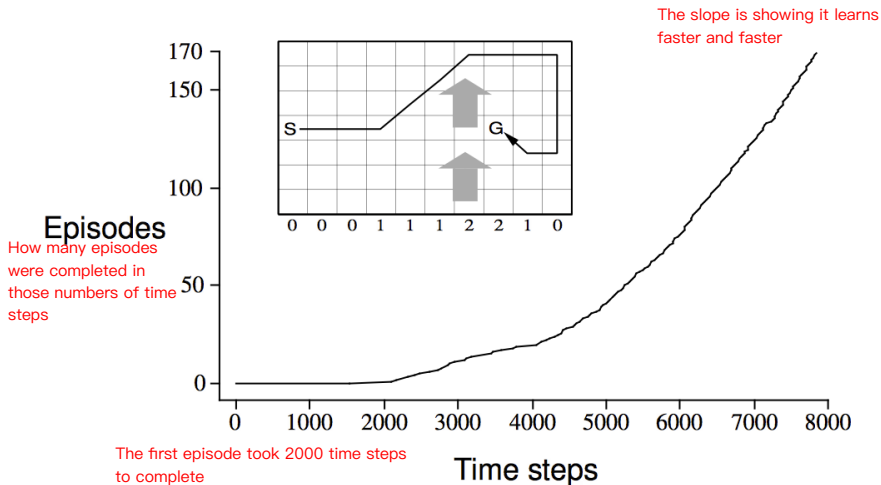


king's  
moves

- Reward = -1 per time-step until reaching goal
- Undiscounted



# Sarsa on the Windy Gridworld



# $n$ -Step Sarsa

- Consider the following  $n$ -step returns for  $n = 1, 2, \infty$ :

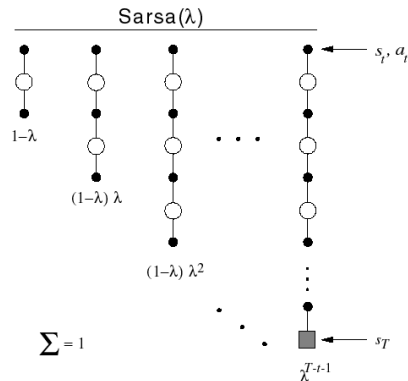
$$\begin{array}{lll}
 n = 1 & (\text{Sarsa}) & q_t^{(1)} = R_{t+1} + \gamma Q(S_{t+1}) \quad \text{1 step ahead} \\
 n = 2 & & q_t^{(2)} = R_{t+1} + \gamma R_{t+2} + \gamma^2 Q(S_{t+2}) \quad \text{2 steps ahead} \\
 & & \vdots \\
 n = \infty & (MC) & q_t^{(\infty)} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{T-1} R_T \quad \text{No bootstrap}
 \end{array}$$

- Define the  $n$ -step Q-return

$$q_t^{(n)} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n Q(S_{t+n})$$

- $n$ -step Sarsa updates  $Q(s, a)$  towards the  $n$ -step Q-return

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha \left( q_t^{(n)} - Q(S_t, A_t) \right)$$

Forward View Sarsa( $\lambda$ )

- The  $q^\lambda$  return combines all  $n$ -step Q-returns  $q_t^{(n)}$
- Using weight  $(1 - \lambda)\lambda^{n-1}$

$$q_t^\lambda = (1 - \lambda) \sum_{n=1}^{\infty} \lambda^{n-1} q_t^{(n)}$$

- Forward-view Sarsa( $\lambda$ )

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha \left( q_t^\lambda - Q(S_t, A_t) \right)$$

# Backward View Sarsa( $\lambda$ )

- Just like TD( $\lambda$ ), we use **eligibility traces** in an online algorithm
- But Sarsa( $\lambda$ ) has one eligibility trace for each state-action pair

$$E_0(s, a) = 0$$

$$E_t(s, a) = \gamma\lambda E_{t-1}(s, a) + \mathbf{1}(S_t = s, A_t = a)$$

- $Q(s, a)$  is updated for every state  $s$  and action  $a$
- In proportion to TD-error  $\delta_t$  and eligibility trace  $E_t(s, a)$

$$\delta_t = R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t)$$

$$Q(s, a) \leftarrow Q(s, a) + \alpha \delta_t E_t(s, a)$$

# Sarsa( $\lambda$ ) Algorithm

Initialize  $Q(s, a)$  arbitrarily, for all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$

Repeat (for each episode):

$E(s, a) = 0$ , for all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$

    Initialize  $S, A$

    Repeat (for each step of episode):

        Take action  $A$ , observe  $R, S'$

        Choose  $A'$  from  $S'$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)

$\delta \leftarrow R + \gamma Q(S', A') - Q(S, A)$

$E(S, A) \leftarrow E(S, A) + \delta$

        For all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$ :

$Q(s, a) \leftarrow Q(s, a) + \alpha \delta E(s, a)$

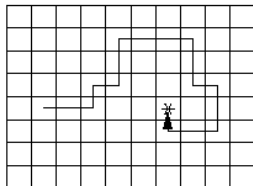
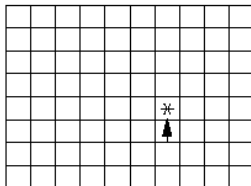
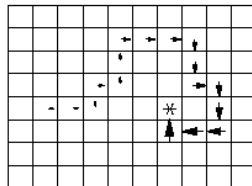
$E(s, a) \leftarrow \gamma \lambda E(s, a)$

$S \leftarrow S'; A \leftarrow A'$

    until  $S$  is terminal

# Sarsa( $\lambda$ ) Gridworld Example

Path taken

Action values increased  
by one-step SarsaAction values increased  
by Sarsa( $\lambda$ ) with  $\lambda=0.9$ 

# Off-Policy Learning

- Evaluate target policy  $\pi(a|s)$  to compute  $v_\pi(s)$  or  $q_\pi(s, a)$
- While following behaviour policy  $\mu(a|s)$

$$\{S_1, A_1, R_2, \dots, S_T\} \sim \mu$$

- Why is this important?
- Learn from observing humans or other agents
- Re-use experience generated from old policies  $\pi_1, \pi_2, \dots, \pi_{t-1}$
- Learn about *optimal* policy while following *exploratory* policy
- Learn about *multiple* policies while following *one* policy

# Importance Sampling

- Estimate the expectation of a different distribution

$$\begin{aligned}\mathbb{E}_{X \sim P}[f(X)] &= \sum P(X)f(X) \\ &= \sum Q(X) \frac{P(X)}{Q(X)} f(X) \\ &= \mathbb{E}_{X \sim Q} \left[ \frac{P(X)}{Q(X)} f(X) \right]\end{aligned}$$



# Importance Sampling for Off-Policy Monte-Carlo

- Use returns generated from  $\mu$  to evaluate  $\pi$
- Weight return  $G_t$  according to similarity between policies
- Multiply importance sampling corrections along whole episode

$$G_t^{\pi/\mu} = \frac{\pi(A_t|S_t)}{\mu(A_t|S_t)} \frac{\pi(A_{t+1}|S_{t+1})}{\mu(A_{t+1}|S_{t+1})} \cdots \frac{\pi(A_T|S_T)}{\mu(A_T|S_T)} G_t$$

- Update value towards *corrected* return

$$V(S_t) \leftarrow V(S_t) + \alpha \left( G_t^{\pi/\mu} - V(S_t) \right)$$

- Cannot use if  $\mu$  is zero when  $\pi$  is non-zero
- Importance sampling can dramatically increase variance

# Importance Sampling for Off-Policy TD

- Use TD targets generated from  $\mu$  to evaluate  $\pi$
- Weight TD target  $R + \gamma V(S')$  by importance sampling
- Only need a single importance sampling correction

$$V(S_t) \leftarrow V(S_t) + \alpha \left( \frac{\pi(A_t|S_t)}{\mu(A_t|S_t)} (R_{t+1} + \gamma V(S_{t+1})) - V(S_t) \right)$$

- Much lower variance than Monte-Carlo importance sampling
- Policies only need to be similar over a single step

# Q-Learning

- We now consider off-policy learning of action-values  $Q(s, a)$
- **No** importance sampling is required
- Next action is chosen using behaviour policy  $A_{t+1} \sim \mu(\cdot|S_t)$
- But we consider alternative successor action  $A' \sim \pi(\cdot|S_t)$
- And update  $Q(S_t, A_t)$  towards value of alternative action

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha (R_{t+1} + \gamma Q(S_{t+1}, A') - Q(S_t, A_t))$$

# Off-Policy Control with Q-Learning

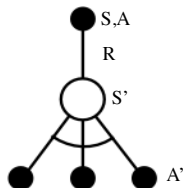
- We now allow both behaviour and target policies to **improve**
- The target policy  $\pi$  is **greedy** w.r.t.  $Q(s, a)$

$$\pi(S_{t+1}) = \operatorname{argmax}_{a'} Q(S_{t+1}, a')$$

- The behaviour policy  $\mu$  is e.g.  **$\epsilon$ -greedy** w.r.t.  $Q(s, a)$
- The Q-learning target then simplifies:

$$\begin{aligned} & R_{t+1} + \gamma Q(S_{t+1}, A') \\ &= R_{t+1} + \gamma Q(S_{t+1}, \operatorname{argmax}_{a'} Q(S_{t+1}, a')) \\ &= R_{t+1} + \max_{a'} \gamma Q(S_{t+1}, a') \end{aligned}$$

# Q-Learning Control Algorithm



$$Q(S, A) \leftarrow Q(S, A) + \alpha \left( R + \gamma \max_{a'} Q(S', a') - Q(S, A) \right)$$

## Theorem

*Q-learning control converges to the optimal action-value function,  $Q(s, a) \rightarrow q_*(s, a)$*

# Q-Learning Algorithm for Off-Policy Control

Initialize  $Q(s, a), \forall s \in \mathcal{S}, a \in \mathcal{A}(s)$ , arbitrarily, and  $Q(\text{terminal-state}, \cdot) = 0$

Repeat (for each episode):

    Initialize  $S$

    Repeat (for each step of episode):

        Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)

        Take action  $A$ , observe  $R, S'$

$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$

$S \leftarrow S'$ ;

    until  $S$  is terminal

# Q-Learning Demo

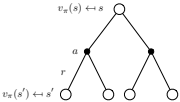
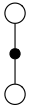
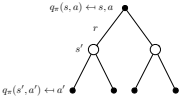
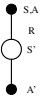
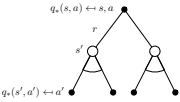
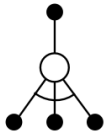
Q-Learning Demo

# Cliff Walking Example





# Relationship Between DP and TD

	<i>Full Backup (DP)</i>	<i>Sample Backup (TD)</i>
Bellman Expectation Equation for $v_{\pi}(s)$	 <p>Iterative Policy Evaluation</p>	 <p>TD Learning</p>
Bellman Expectation Equation for $q_{\pi}(s, a)$	 <p>Q-Policy Iteration</p>	 <p>Sarsa</p>
Bellman Optimality Equation for $q_{*}(s, a)$	 <p>Q-Value Iteration</p>	 <p>Q-Learning</p>

# Relationship Between DP and TD (2)

<i>Full Backup (DP)</i>	<i>Sample Backup (TD)</i>
Iterative Policy Evaluation $V(s) \leftarrow \mathbb{E}[R + \gamma V(S') \mid s]$	TD Learning $V(S) \stackrel{\alpha}{\leftarrow} R + \gamma V(S')$
Q-Policy Iteration $Q(s, a) \leftarrow \mathbb{E}[R + \gamma Q(S', A') \mid s, a]$	Sarsa $Q(S, A) \stackrel{\alpha}{\leftarrow} R + \gamma Q(S', A')$
Q-Value Iteration $Q(s, a) \leftarrow \mathbb{E}\left[R + \gamma \max_{a' \in \mathcal{A}} Q(S', a') \mid s, a\right]$	Q-Learning $Q(S, A) \stackrel{\alpha}{\leftarrow} R + \gamma \max_{a' \in \mathcal{A}} Q(S', a')$

where  $x \stackrel{\alpha}{\leftarrow} y \equiv x \leftarrow x + \alpha(y - x)$

# Questions?